

Exhibit 1

Overview of the TREC 2011 Legal Track

Notebook Draft 2011.10.24

Maura R. Grossman
Gordon V. Cormack
Bruce Hedin
Douglas W. Oard

Abstract

The TREC 2011 Legal Track consisted of one task: the *learning task*, which captured elements of the TREC 2010 learning and interactive tasks. Participants were required to rank the entire corpus of 670,000 documents by their estimate of the probability of relevance to each of 3 topics, and also to give a quantitative estimate of that probability. Participants were permitted to request up to 1,000 *relevance determinations* from a *Topic Authority* for each topic. Participants elected either to use only these relevance determinations in preparing *automatic* submissions, or to augment these determinations with their own manual review in preparing *technology assisted* submissions. We provide a brief overview of the task, and preliminary results as of October 24, 2011. More detailed results will be available to TREC participants during the conference at the web address <http://plg1.uwaterloo.ca/trec11-assess> .

1 Introduction

We are concerned with the document selection and review component of the *e-discovery* process, for which the objective is to identify as nearly as practicable all documents from a collection that are responsive to a *request for production* in civil litigation, while minimizing the number of unresponsive documents that are identified by the method.

The learning task models the scenario in which a senior attorney – the *Topic Authority* – is charged with interpreting the request for production, communicating that interpretation to a review team, and producing responsive documents to the requesting party. TREC participants play the role of the review team.

At the outset, the Topic Authority reads the request and prepares a set of coding guidelines. The request and the guidelines are given to participants, and an initial *kick-off call* allows interested participants to ask the Topic Authority how to interpret the request for production.

Over the course of several weeks, each participant is entitled to request feedback from the Topic Authority on a number of documents from the collection. This feedback consists of a simple binary *relevance determination*: Participants are informed whether the Topic Authority determines each document to be responsive or not. No other communication with the Topic Authority is permitted.

Teams from ten different groups participated in the Legal Track; the names of the teams, as well as the prefix used to label the team's results, are shown in table 1.

2 Document Collection

The document collection for the TREC 2011 Legal Track is identical to that used for TREC 2010. It was derived from the EDRM Enron Dataset, version 2, prepared by ZL Technologies in consultation with the Legal Track coordinators, and hosted by EDRM. The EDRM dataset consists of 1.3 million email messages captured by the Federal Energy Review Commission (FERC) from Enron, in the course of its investigation of Enron's collapse. ZL acquired the dataset from Loughheed Systems (formerly Aspen Systems) who captured

Participating Organization	Run Prefix
Beijing University of Posts and Telecommunications	pri
Helioid	HEL
Indian Statistical Institute	ISI
OpenText	ot
Recommind	rec
TCDI	TCD
University of Melbourne	mlb
University of South Florida	USF
University of Waterloo	UW
Ursinus College	URS

Table 1: Organizations participating in the TREC 2011 Legal Track.

and maintain the dataset on behalf of FERC. The EDRM dataset is available in two formats: EDRM XML and PST. The EDRM XML version contains a text rendering of each email message and attachment, as well as the original native format. The PST version contains the same messages, in a Microsoft proprietary format used by many commercial tools.

Both versions of the dataset approach 100GB in size, presenting an obstacle to participants. Furthermore, there are a large number of duplicate email messages in the dataset, that were captured more than once by Loughheed. For TREC, a list of 470,000 distinct messages were identified as canonical; all other messages duplicate one of the canonical messages. These messages contain about 200,000 attachment files; together these 470,000 messages plus 200,000 attachments form the 670,000 documents of the TREC 2010/2011 Legal Track collection. Text and native versions of these documents were made available to participants, along with a mapping from the EDRM XML and PST files to their canonical counterparts in the TREC collection.

3 Relevance Assessments

In order to measure the efficacy of TREC participant efforts in the two tasks, it is necessary to compare their results to a *gold standard* indicating whether or not each document in the collection is relevant to a particular discovery request. The learning task used three distinct discovery requests. Ideally, a gold standard would indicate the relevance of each document to each topic.

It is impractical to use human assessors to render these two million assessments. Instead, a sample of documents was identified for each topic, and assessors were asked to code only the documents in the sample as relevant or not. The evaluation was conducted in two phases: preliminary and final. At the time of writing, only the preliminary evaluation was complete.

3.1 Preliminary Evaluation

A total of 16,999 documents – about 5,600 per topic – were selected and assessed to form the preliminary gold standard. The documents were selected according to four criteria:

1. All documents that were identified by the coordinators, in the course of composing the topics before the start of the task, to be potentially relevant;
2. All documents submitted by any team for relevance determination;
3. All documents ranked among the 100 most probably relevant by any submission;
4. A uniform random sample of the remaining documents.

11,612 documents (the *100 stratum*) were selected according to one or more of the first three criteria; 5,387 documents (the *1000 stratum*) were sampled according to the fourth. All documents in the 100 stratum were

assessed, regardless of whether or not a relevance determination had been previously rendered by the Topic Authority. Each document in the 1000 stratum was given to two assessors; that is, each sampled document was assessed twice.

The learning task assessments were rendered by four professional review companies, who volunteered their services. Three of the companies used a Web-based platform developed by the coordinators to view scanned documents and to record their relevance judgments. To avoid problems with local rendering software on each assessor's workstation, the assessors made their judgments based on pdf-formatted versions of the documents, as opposed to the original native format documents. The fourth review company downloaded the pdf documents and conducted the review on their own platform.

Assessors were provided with orientation and detailed guidelines created by a Topic Authority. The review platform included a "seek assistance" link which assessors were encouraged to use to request that the Topic Authority resolve any uncertainty that may have arisen. Assessors were instructed to make a relevance judgment of relevant (R), not relevant (N), or broken (B) for every document in their bins. The latter code reflects the fact that a small percentage of documents from the EDRM dataset are malformed and therefore cannot be assessed.

Once the preliminary assessments were complete, quality assurance was conducted by having the Topic Authority adjudicate conflicting assessments, which occurred in two cases:

1. For documents selected according to criterion 2 above, the Topic Authority's initial relevance determination and the assessor's relevance judgment differed; or
2. For documents selected according to criterion 4 above, i.e. the 1000 stratum, the two assessors' judgments differed.

The Topic Authority adjudicated all conflicting documents together, with no indication of which documents had been subject to a previous relevance determination, or what that determination had been.

The preliminary gold standard consists of

- The assessor's judgment, for documents without conflicting assessments; and,
- The Topic Authority's final judgment, for documents with conflicting assessments.

The preliminary gold standard, along with the toolkit used for the preliminary evaluation, may be found on the web: <http://durum0.uwaterloo.ca/trec/legal10-results> .

3.2 Final Evaluation

At the time of writing, additional assessments and quality assurance measures are being undertaken to improve the accuracy and reliability of the evaluation results. The set of documents identified by any submission as one of the 5,000 most probably relevant – the *5000 stratum* – has been sampled for assessment by the professional review companies. For topic 402, a second sample of the 1000 stratum is also being assessed. Additional redundant assessments are being conducted, and the toolkit is being enhanced to adjust the estimated evaluation measures to compensate for random errors in relevance assessments.

Final results will be available to participants at the TREC workshop in November and, at the same time, on the web: <http://plg1.uwaterloo.ca/trec11-assess> .

4 The Task

The learning task models the use of automated or semi-automated methods to guide review strategy for a multi-stage document review effort, organized as follows:

1. **Initial search and assessment.** The responding party analyzes the production request. Using ad hoc methods the team identifies a *seed set* of potentially responsive documents, and assesses each as responsive or not.

2. **Learning by example.** A learning method is used to rank the documents in the collection from most to least likely to be responsive to the production request, and to estimate this likelihood for each document. The input to the learning method consists of the seed set, the assessments for the seed set, and the unranked collection; the output is a ranked list consisting of the document identifier and a probability of responsiveness for each document in the collection.

The two components of learning by example – ranking and estimation – may be accomplished by the same method or by different methods. Either may be automated or manual. For example, ranking may be done using an information retrieval method or by human review using a five-point scale. Estimation may be done in the course of ranking or, for example, by sampling and reviewing documents at representative ranks.

3. **Review process.** A review process may be conducted, with strategy guided by the ranked list. One possible strategy is to review documents in order, thus discovering as many responsive documents as possible for a given amount of effort. Another possible strategy is triage: to review only mid-ranked documents, deeming, without further review, the top-ranked ones to be responsive, and the bottom-ranked ones to be non-responsive.

Review strategy may be guided not only by the order of the ranked list, as outlined above, but also by the estimated effectiveness of various alternatives. Consider the strategy of reviewing the top-ranked documents. Where should a *cut* be made so that documents above the cut are reviewed and documents below are not? For triage, where should the two necessary cuts be made?

Practically every review strategy decision boils down to the question,

Of this particular set of documents, how many are responsive and how many are not?

This question itself can be reduced to,

What is the probability of each document in the set being relevant?

Given an answer to the second question, the answer to the first is simply the sum of the probabilities. For this reason, participants in the learning track were required to provide an estimate of the probability of relevance for each document in the collection. Using these estimates the documents were ranked from most likely to least likely relevant. At each rank, the estimated number of relevant documents – the sum of the probabilities up to that rank – was computed, and used to estimate recall, precision and F_1 .

4.1 Submission Phases

For each topic, teams were required to submit an *initial* set of probability estimates prior to requesting any relevance determinations from the Topic Authority. Thereafter, teams were required to submit *interim* sets of probability estimates in order to receive more than 100 and more than 300 relevance determinations. A team was allowed to request at most 1,000 relevance determinations per topic.

Each team was required to submit a *final* set of probability estimates once it had received all the relevance determinations requested by the team.

In a final *mopup* phase, all relevance determinations requested by all teams were distributed to all teams, who had the opportunity to submit a *mopup* set of probability estimates.

In this overview, we report results separately for the *initial*, *final* and *mopup* submissions. The run identifiers for the various phases may be distinguished by their final symbol: initial submissions end in “I”; final submissions end in “F”; and mopup submissions end in “M”.

4.2 Participation Categories

Participants were asked to declare each run to be *automatic* or *technology assisted*. Automatic runs were allowed to use manual query formulation, but human review of the document collection (other than that

Topic	Number of Responsive Documents
401	30,853
402	1,920
403	1,239

Table 2: Estimated number of responsive documents for each topic.

provided by TREC via responsiveness determinations) was not allowed. Technology assisted runs were allowed to avail themselves of any amount of human review. Participants were asked to state the number of documents reviewed, as well as the number of hours spent – both reviewing documents and configuring the system. The participation category of a run is specified by the penultimate character in its name: “A” for automatic; and “T” for technology assisted. For example, the run named gggxxxAF is a final run, automatic participation, by the group whose run prefix is ggg.

4.3 Topics

There were three topics: 401, 402 and 403.

- Topic 401 (Topic Authority: Kevin F. Brady, Connolly Bove Lodge & Hutz LLP.)
All documents or communications that describe, discuss, refer to, report on, or relate to the design, development, operation, or marketing of enrononline, or any other online service offered, provided, or used by the Company (or any of its subsidiaries, predecessors, or successors-in-interest), for the purchase, sale, trading, or exchange of financial or other instruments or products, including but not limited to, derivative instruments, commodities, futures, and swaps.
- Topic 402 (Topic Authority: Brendan M. Schulman, Kramer Levin Naftalis & Frankel LLP.)
All documents or communications that describe, discuss, refer to, report on, or relate to whether the purchase, sale, trading, or exchange of over-the-counter derivatives, or any other actual or contemplated financial instruments or products, is, was, would be, or will be legal or illegal, or permitted or prohibited, under any existing or proposed rule(s), regulation(s), law(s), standard(s), or other proscription(s), whether domestic or foreign.
- Topic 403 (Topic Authority: Robert Singleton, Squire, Sanders & Dempsey (US) LLP.)
All documents or communications that describe, discuss, refer to, report on, or relate to the environmental impact of any activity or activities undertaken by the Company, including but not limited to, any measures taken to conform to, comply with, avoid, circumvent, or influence any existing or proposed rule(s), regulation(s), law(s), standard(s), or other proscription(s), such as those governing environmental emissions, spills, pollution, noise, and/or animal habitats.

5 Preliminary Results

Based on the preliminary gold standard, we estimate the total number of responsive documents in the collection to be as shown in Table 2.

The following sections detail the results achieved by the various submissions. We note that not all teams undertook every topic, and not all teams submitted initial or mopup submissions. We note further that some recall estimates in the summary statistics exceed 1.000, due to random error in the preliminary evaluation.

5.1 Gain Curves

A *gain curve* shows the fraction of responsive documents (*recall*) as a function of the number of documents produced, when the documents are produced in order from most likely to least likely relevance. Thus each submission has an associated gain curve. Figure 1 shows the gain curves for the Topic 401 initial

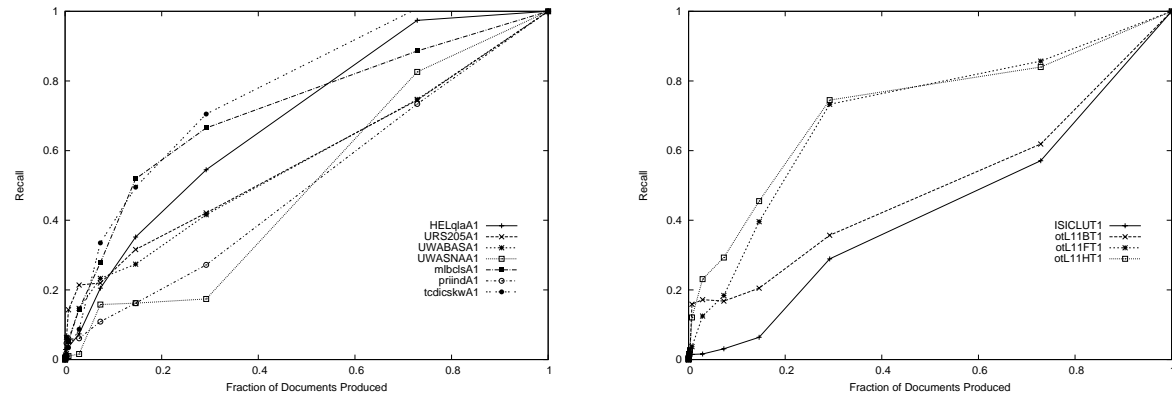


Figure 1: Gain curves for Topic 401 initial submissions (Left: automatic; Right: technology-assisted).

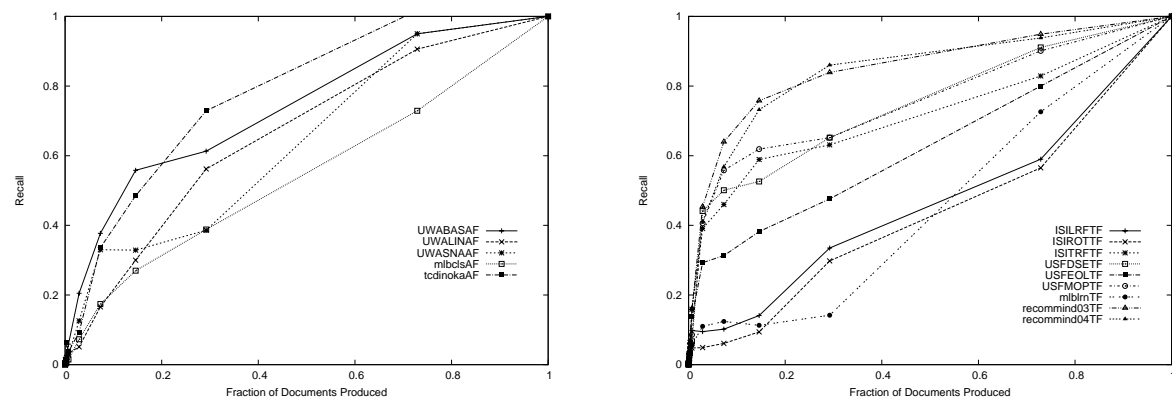


Figure 2: Gain curves for Topic 401 final submissions (Left: automatic; Right: technology-assisted).

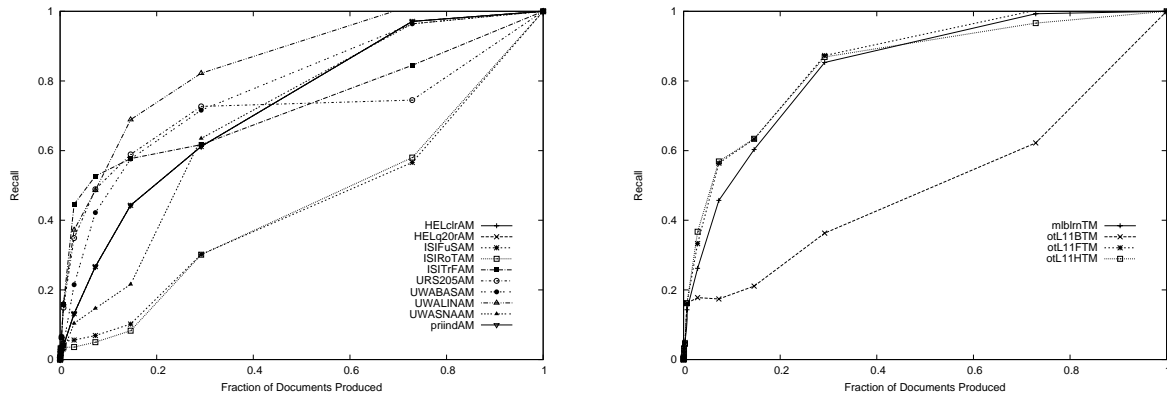


Figure 3: Gain curves for Topic 401 mopup submissions (Left: automatic; Right: technology-assisted).

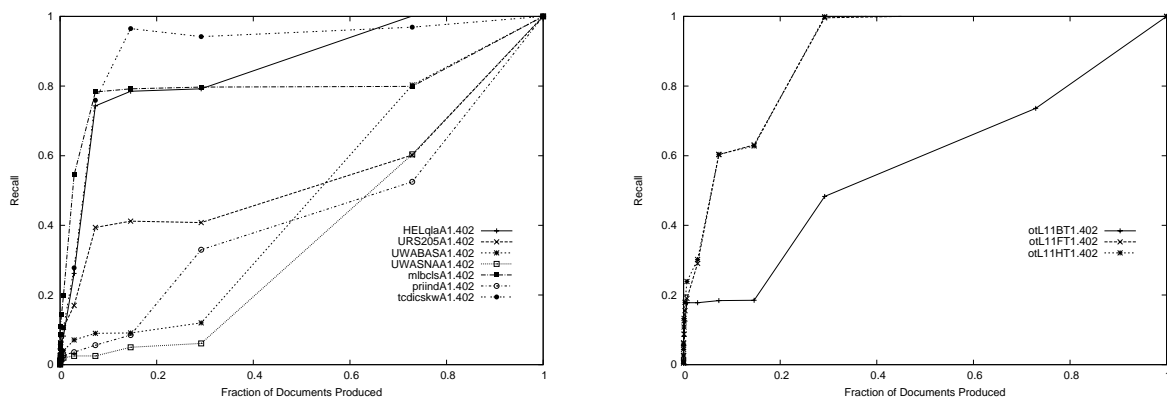


Figure 4: Gain curves for Topic 402 initial submissions (Left: automatic; Right: technology-assisted).

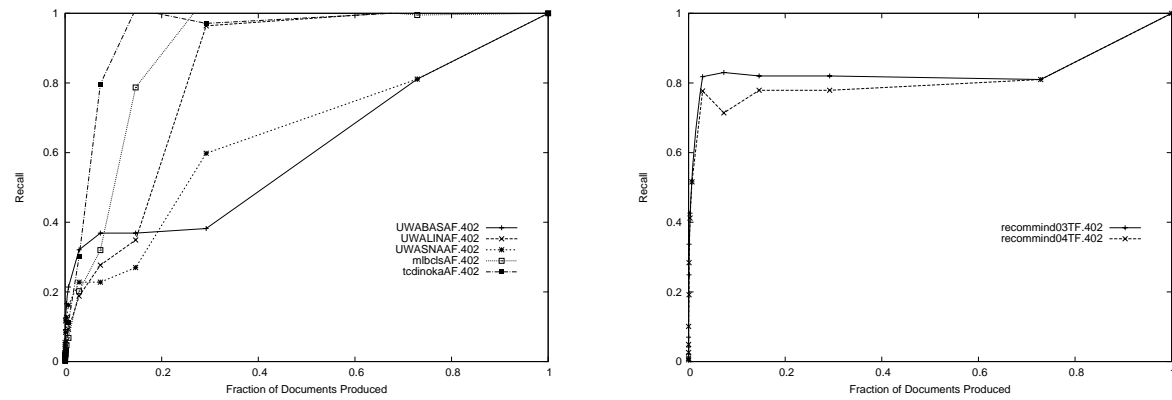


Figure 5: Gain curves for Topic 402 final submissions (Left: automatic; Right: technology-assisted).

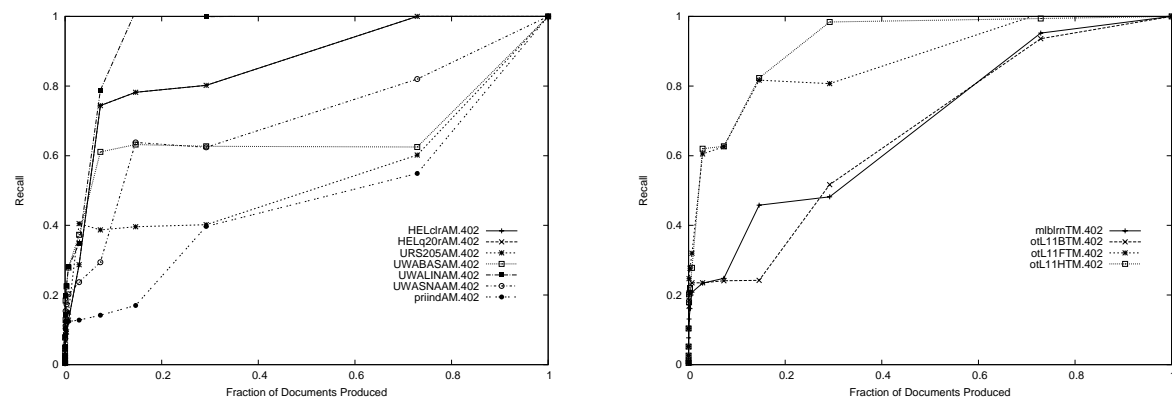


Figure 6: Gain curves for Topic 402 mopup submissions (Left: automatic; Right: technology-assisted).

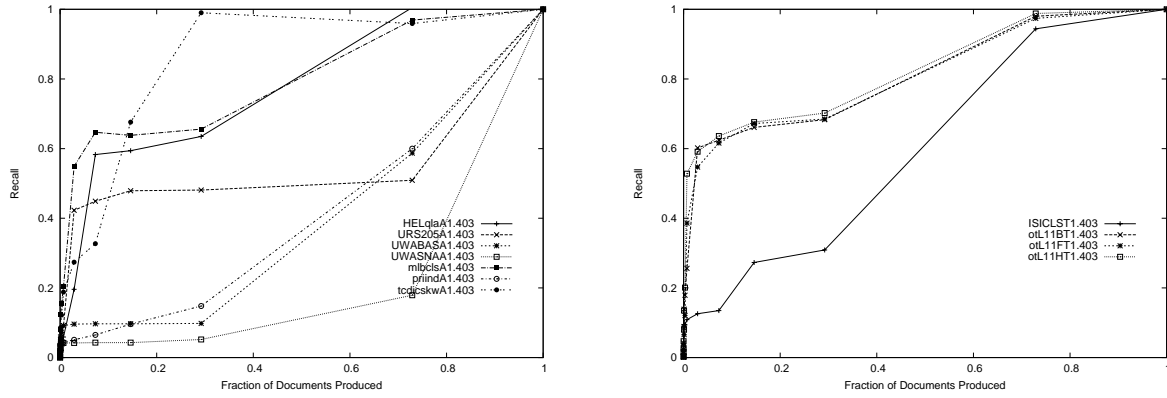


Figure 7: Gain curves for Topic 403 initial submissions (Left: automatic; Right: technology-assisted).

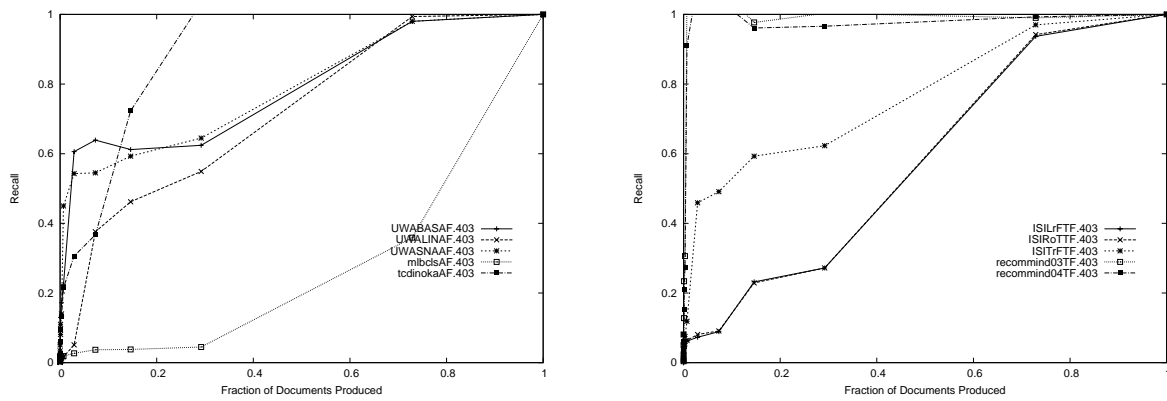


Figure 8: Gain curves for Topic 403 final submissions (Left: automatic; Right: technology-assisted).

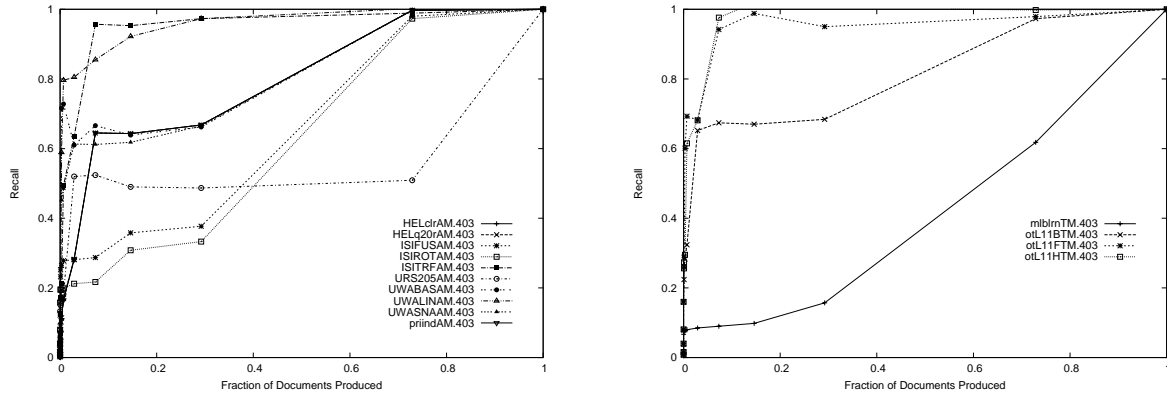


Figure 9: Gain curves for Topic 403 mopup submissions (Left: automatic; Right: technology-assisted).

submissions. The automatic submissions are shown on the left and the technology-assisted submissions on the right. Figures 2 and 3 show the corresponding curves for the final and mopup submissions.

Figures 4 through 6 show the corresponding information for Topic 402; Figures 7 through 9 for Topic 403.

5.2 Precision, Recall and F_1

Because participants submitted a probability estimate rather than a set of responsive documents, it is necessary to specify a *cutoff rank*, c , such that the c documents with the highest probability of relevance are deemed to be responsive, and the remainder are deemed to be non-responsive. Once the cutoff rank is chosen, recall, precision and F_1 may be calculated in the normal way.

For this set of evaluation measures we assume that the cutoff rank is to be chosen so as to maximize F_1 . Once the gold standard has been constructed, it is a simple matter to try all possible values of c and to choose the one that yields the maximum value of F_1 . We call this value the *hypothetical* F_1 , because it could be achieved, but only if the appropriate value of c were somehow determined without the aid of the gold standard.

Since c must be determined without knowledge of the gold standard, we must rely on the probability estimates contained in the submissions. As noted above, the probability estimates for the top-ranked c documents may be summed to yield an estimate of the number of responsive documents were the cut to be made at c . From this estimate we may easily derive estimates of recall, precision and F_1 , and select the value of c that yields the largest F_1 estimate. We call this value the *actual* F_1 , as it can be achieved using only the information contained in the submission.

Tables 3 through 11 show the hypothetical and actual F_1 values, along with the recall and precision values at the cutoff rank that achieves the corresponding F_1 value. Separate tables are shown for initial, final and mopup runs, and for each topic.

6 Discussion

It is apparent from the gain curves that the best systems are able to identify the vast majority of responsive documents with a cutoff value that includes only a tiny fraction of the collection. For Topic 401, this level of recall is achieved with a cutoff of about 10%. For topics 402 and 403, the number is much lower.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
otL11HT1	0.372	0.276	0.572	0.302	0.197	0.646
URS205A1	0.343	0.223	0.745	0.168	0.092	0.987
otL11BT1	0.324	0.199	0.868	0.284	0.177	0.727
tclicskwA1	0.279	0.436	0.205	0.115	1.005	0.061
mlbclsA1	0.259	0.423	0.187	0.095	0.971	0.050
otL11FT1	0.213	0.746	0.124	0.148	0.167	0.133
UWABASA1	0.198	0.150	0.291	0.089	0.692	0.048
HELqlaA1	0.183	0.349	0.124	0.086	1.000	0.045
priindA1	0.127	0.070	0.724	0.111	0.061	0.600
UWASNA1	0.122	0.168	0.096	0.091	0.689	0.049
ISICLUT1	0.088	0.990	0.046	0.000	0.000	0.875

Table 3: Topic 401 initial submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
recommind03TF	0.588	0.585	0.591	0.476	0.319	0.939
recommind04TF	0.581	0.435	0.874	0.510	0.350	0.937
USFDSETF	0.560	0.486	0.660	0.128	0.077	0.370
USFMOPTF	0.543	0.481	0.623	0.176	0.113	0.406
ISITRFTF	0.505	0.433	0.606	0.000	0.000	0.875
USFEOLTF	0.426	0.322	0.631	0.376	0.242	0.847
UWABASAF	0.339	0.510	0.254	0.109	0.983	0.058
tcclnokaAF	0.286	0.439	0.212	0.088	0.995	0.046
UWASNAAF	0.253	0.331	0.205	0.109	0.983	0.058
mlblnTF	0.179	0.111	0.457	0.040	0.021	0.322
ISILRFTF	0.176	0.103	0.602	0.000	0.000	0.125
UWALINAF	0.160	0.273	0.114	0.115	0.831	0.062
mlbclsAF	0.143	0.161	0.129	0.020	0.011	0.173
ISIROTF	0.122	0.067	0.657	0.000	0.000	0.875

Table 4: Topic 401 final submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
ISITrFAM	0.578	0.456	0.788	0.001	0.000	1.000
otL11HTM	0.505	0.397	0.694	0.332	0.200	0.997
otL11FTM	0.484	0.368	0.708	0.084	0.044	0.834
UWALINAM	0.458	0.401	0.533	0.272	0.796	0.164
URS205AM	0.446	0.409	0.491	0.427	0.378	0.492
mlblnTM	0.377	0.359	0.398	0.376	0.355	0.398
UWABASAM	0.373	0.422	0.335	0.109	0.995	0.058
otL11BTM	0.337	0.208	0.893	0.318	0.193	0.896
HELq20rAM	0.246	0.451	0.169	0.142	0.950	0.077
HELclrAM	0.246	0.451	0.169	0.142	0.950	0.077
UWASNAAM	0.177	0.699	0.101	0.109	0.995	0.058
ISIFuSAM	0.141	0.076	0.992	0.001	0.000	1.000
priindAM	0.095	0.148	0.070	0.062	0.034	0.400
ISIRoTAM	0.088	0.989	0.046	0.001	0.000	1.000

Table 5: Topic 401 mopup submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
otL11HT1	0.185	0.131	0.314	0.065	0.289	0.036
otL11FT1	0.168	0.130	0.238	0.047	0.308	0.025
mlbclsA1	0.156	0.139	0.177	0.006	0.992	0.003
otL11BT1	0.130	0.091	0.232	0.120	0.138	0.106
URS205A1	0.086	0.057	0.176	0.036	0.139	0.021
tclicskwA1	0.078	0.736	0.041	0.007	0.970	0.004
HELqlaA1	0.078	0.052	0.155	0.006	1.000	0.003
UWABASA1	0.034	0.035	0.034	0.005	0.538	0.002
UWASNAA1	0.025	0.015	0.092	0.003	0.346	0.001
priindA1	0.023	0.015	0.053	0.022	0.026	0.019

Table 6: Topic 402 initial submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
recommind03TF	0.588	0.423	0.960	0.477	0.384	0.629
recommind04TF	0.455	0.440	0.470	0.395	0.387	0.404
UWABASAF	0.172	0.156	0.191	0.007	1.005	0.004
UWASNAAF	0.154	0.117	0.225	0.007	1.004	0.004
tcclnokaAF	0.087	0.776	0.046	0.006	0.999	0.003
UWALINAF	0.059	0.067	0.053	0.020	0.752	0.010
mlbclsAF	0.046	0.047	0.046	0.037	0.186	0.020

Table 7: Topic 402 final submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
otL11FTM	0.345	0.234	0.655	0.224	0.126	1.000
otL11HTM	0.289	0.192	0.581	0.286	0.190	0.580
otL11BTM	0.273	0.166	0.781	0.209	0.199	0.220
UWALINAM	0.262	0.201	0.377	0.091	0.746	0.048
UWABASAM	0.250	0.199	0.337	0.007	1.007	0.004
priindAM	0.221	0.124	1.000	0.082	0.124	0.061
UWASNAM	0.209	0.162	0.293	0.007	1.007	0.004
mlblmTM	0.175	0.142	0.226	0.123	0.066	0.858
URS205AM	0.158	0.100	0.374	0.068	0.395	0.037
HELclrAM	0.111	0.073	0.226	0.010	1.001	0.005
HELq20rAM	0.110	0.074	0.213	0.010	1.002	0.005

Table 8: Topic 402 mopup submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
otL11HT1	0.335	0.580	0.235	0.106	0.539	0.059
otL11FT1	0.234	0.339	0.179	0.078	0.494	0.042
otL11BT1	0.188	0.552	0.114	0.104	0.587	0.057
mlbclsA1	0.149	0.434	0.090	0.004	0.986	0.002
ISICLST1	0.110	0.078	0.183	0.008	0.004	0.625
tcdictskwA1	0.098	0.149	0.073	0.005	0.958	0.002
URS205A1	0.083	0.352	0.047	0.033	0.105	0.019
UWABASA1	0.066	0.057	0.080	0.001	0.128	0.000
priindA1	0.063	0.036	0.234	0.019	0.044	0.012
UWASNAA1	0.035	0.040	0.031	0.001	0.122	0.000
HELqlaA1	0.032	0.578	0.016	0.004	1.000	0.002

Table 9: Topic 403 initial submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
recommind03TF	0.623	1.332	0.407	0.247	0.258	0.237
recommind04TF	0.574	0.903	0.421	0.065	0.996	0.033
UWASNAAF	0.241	0.491	0.160	0.005	0.984	0.002
UWABASAF	0.146	0.127	0.171	0.005	0.984	0.002
ISITrFTF	0.139	0.443	0.083	0.008	0.004	0.625
tcclinokaAF	0.107	0.166	0.079	0.004	0.996	0.002
ISIRoTTF	0.090	0.048	0.787	0.008	0.004	0.625
ISILrFTF	0.046	0.032	0.078	0.000	0.000	0.000
UWALINAF	0.035	0.322	0.019	0.009	0.512	0.005
mlbclsAF	0.031	0.017	0.221	0.002	0.367	0.001

Table 10: Topic 403 final submission results.

run	Hypothetical			Actual		
	F_1	Recall	Precision	F_1	Recall	Precision
UWASNAAM	0.720	0.702	0.739	0.005	0.997	0.002
otL11FTM	0.612	0.664	0.567	0.345	0.215	0.875
UWABASAM	0.578	0.827	0.444	0.005	0.998	0.002
UWALINAM	0.467	0.660	0.362	0.048	0.820	0.025
otL11HTM	0.376	0.249	0.759	0.310	0.270	0.363
ISIFUSAM	0.340	0.228	0.670	0.013	0.006	1.000
otL11BTM	0.337	0.210	0.850	0.228	0.281	0.192
priindAM	0.325	0.195	0.984	0.123	0.195	0.090
ISIROTAM	0.325	0.195	0.984	0.013	0.006	1.000
ISITRFAM	0.249	0.508	0.165	0.013	0.006	1.000
URS205AM	0.239	0.145	0.672	0.043	0.521	0.022
mlblmTM	0.134	0.073	0.734	0.055	0.028	0.795
HELq20rAM	0.093	0.097	0.090	0.007	1.004	0.003
HELclrAM	0.090	0.090	0.089	0.007	1.004	0.003

Table 11: Topic 403 mopup submission results.

The evaluation of recall, precision and F_1 shows overwhelmingly that the submitted probability estimates are very poor. If the probability estimates were accurate, they would yield an actual F_1 value close to the hypothetical F_1 . Estimation is no mere academic exercise. A fundamental question in e-discovery is: when have I done enough? Estimation is essential to answering this question. The coordinators hope that these results will provide impetus to discover better estimation techniques.

As we have noted, these results are preliminary. It may be that errors in the gold standard cause system performances to be underestimated or (perhaps less likely) to be overestimated. The results contained here will be superceded by final results in which we will seek to mitigate these errors.

7 Conclusion

This is the sixth year of the TREC Legal Track, and our third year of building test collections based on Enron email [1, 5, 4, 3, 2]. Relevance judgments are now available for 14 topical production requests in addition to the (2010) review for privilege. The Legal Track will continue in 2012, when we anticipate having a new collection available for use by participants. We look forward to discussing what we have learned this year and the opportunities for 2012 when we meet this November in Gaithersburg.

References

- [1] J. Baron, D. Lewis, and D. Oard. Trec 2006 legal track overview. In *Proc. 15th Text REtrieval Conference*, 2006.
- [2] G. Cormack, M. Grossman, B. Hedin, and D. Oard. Overview of the trec 2010 legal track. In *Proc. 19th Text REtrieval Conference*, 2010. to appear.
- [3] B. Hedin, D. Oard, S. Tomlinson, and J. Baron. Overview of the trec 2009 legal track. In *Proc. 18th Text REtrieval Conference*, 2009.
- [4] D. Oard, B. Hedin, S. Tomlinson, and J. Baron. Overview of the trec 2008 legal track. In *Proc. 17th Text REtrieval Conference*, 2008.
- [5] S. Tomlinson, D. Oard, J. Baron, and P. Thompson. Overview of the trec 2007 legal track. In *Proc. 16th Text REtrieval Conference*, 2008.